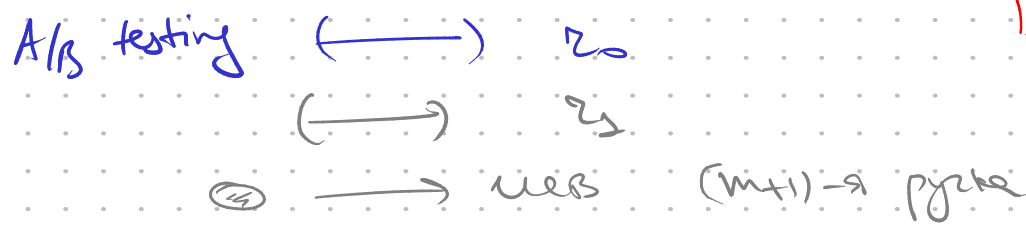
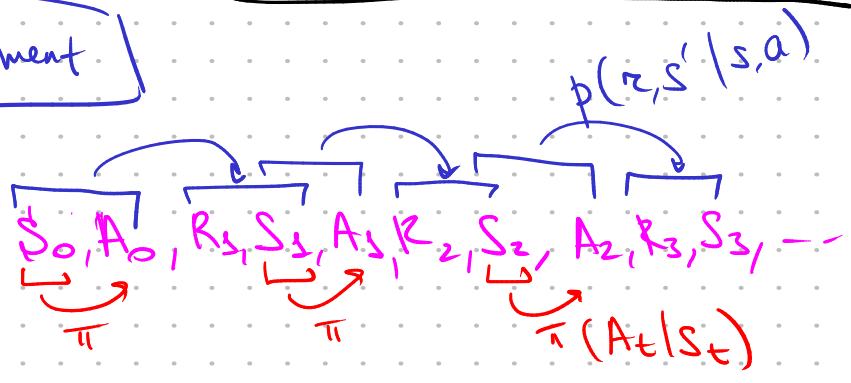
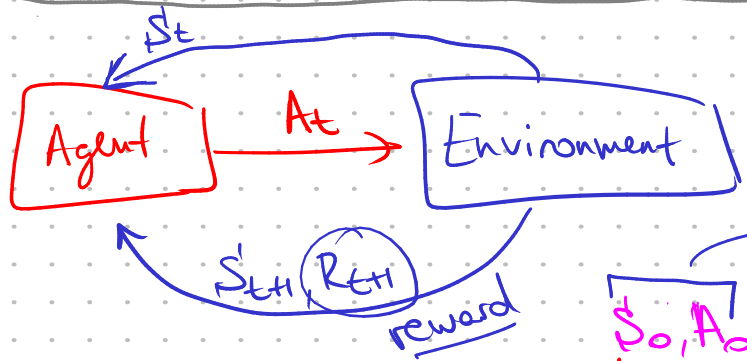


- a_1 c_1 $r_1(a_1, c_1)$
- a_2 c_2 $r_2(a_2, c_2)$
- \vdots \vdots \vdots
- a_t c_t $r_t(a_t, c_t)$
- \vdots \vdots \vdots



MDP - Markov Decision Process



$p(r, s' | s, a)$, $p(s' | s, a)$,
 $r(s, a) = E[R_{t+1} | S_t = s, A_t = a]$
 $r(s, a, s') = E[R_{t+1} | S_t = s, A_t = a, S_{t+1} = s']$

① Что означают? G_t - expected return

Episodic tasks: $G_t = R_{t+1} + R_{t+2} + \dots + R_T$

Continuous tasks: $G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$

$G_t = R_{t+1} + \gamma \cdot G_{t+1}$

$\gamma = 1, R_{T+k} = 0$

$\sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$
 $\leq \max R \cdot \frac{1}{1-\gamma}$

② Value functions

- State value function

$V_{\pi}(s) = E_{\pi} [G_t | S_t = s] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$
 $A_t = \pi(S_t)$

- action value function

$Q_{\pi}(s, a) = E_{\pi} [G_t | S_t = s, A_t = a] = E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, A_t = a \right]$

$$\pi^*(s) = \operatorname{argmax}_{\pi} V_{\pi}(s) = \operatorname{argmax}_{\pi, a} Q_{\pi}(s, a)$$

$$Q_{\pi}(s, a) = \mathbb{E}_{R_{t+1}, S_{t+1}, A_{t+1}, S_{t+2}, \dots} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s, A_t = a] =$$

$$= \mathbb{E}_{\substack{R_{t+1}, S_{t+1}}} [\mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1})]] =$$

$$= \sum_{r, s'} p(r, s' | s, a) \left(r + \mathbb{E}_{\pi} [\gamma R_{t+2} + \gamma^2 R_{t+3} + \dots \mid S_t = s, A_t = a, S_{t+1} = s'] \right)$$

$$Q_{\pi}(s, a) = \sum_{r, s'} p(r, s' | s, a) (r + \gamma V_{\pi}(s'))$$

$$V_{\pi}(s) = \mathbb{E}_{A_t, R_{t+1}, S_{t+1}, \dots} [G_t \mid S_t = s] = \sum_a \pi(a | s) \mathbb{E}_{\pi} [-]$$

$$V_{\pi}(s) = \sum_a \pi(a | s) \cdot Q_{\pi}(s, a)$$

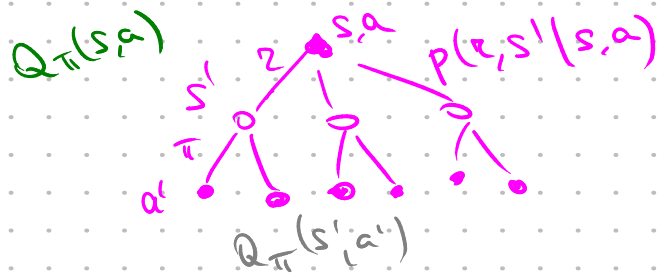
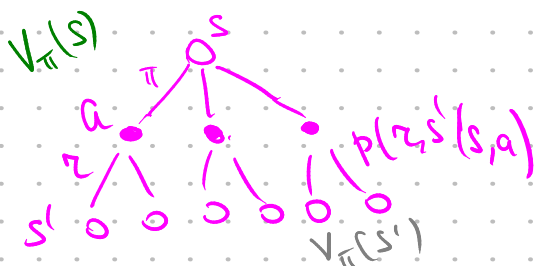
$$V_{\pi}(s) = \sum_a \pi(a | s) \sum_{r, s'} p(r, s' | s, a) (r + \gamma V_{\pi}(s'))$$

Bellman equation

$$\begin{pmatrix} V_{\pi}(s_1) \\ \vdots \\ V_{\pi}(s_N) \end{pmatrix} = \begin{pmatrix} \text{---} \\ \text{---} \\ \text{---} \end{pmatrix} \begin{pmatrix} V_{\pi}(s_1) \\ \vdots \\ V_{\pi}(s_N) \end{pmatrix} \quad \pi = f(v)$$

$$Q_{\pi}(s, a) = \sum_{r, s'} p(r, s' | s, a) \left(r + \gamma \sum_{a'} \pi(a' | s') Q_{\pi}(s', a') \right)$$

backup diagrams



$$V_*(s) = \max_{\pi} V_{\pi}(s) = \max_{\pi} \left(\sum_a \pi(a|s) \sum_{z,s'} p(z,s'|s,a) (r + \gamma V_{\pi}(s')) \right)$$

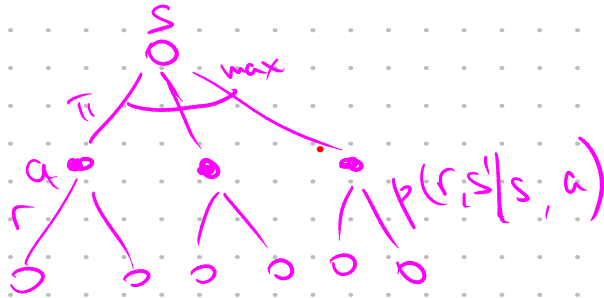
$$= \max_{\pi(a|s)} \max_{\pi} \left(\sum_a \pi(a|s) \sum_{z,s'} p(z,s'|s,a) (r + \gamma V_{\pi}(s')) \right)$$

$$\geq \max_a \sum_{z,s'} p(z,s'|s,a) (r + \gamma \max_{\pi} V_{\pi}(s'))$$

$$V_*(s) = \max_a \left[\sum_{z,s'} p(z,s'|s,a) (r + \gamma V_*(s')) \right]$$

$\bar{x} = f(\bar{x})$

Bellman equation



$$Q_*(s,a) = \max_{\pi} Q_{\pi}(s,a) = \max_{\pi} \sum_{z,s'} p(z,s'|s,a) (r + \gamma \sum_{a'} \pi(a'|s') Q_{\pi}(s',a'))$$

$$= \sum_{z,s'} p(z,s'|s,a) (r + \gamma \max_{\pi(a|s')} \max_{\pi} \left[\sum_{a'} \pi(a'|s') Q_{\pi}(s',a') \right])$$

$$Q_*(s,a) = \sum_{z,s'} p(z,s'|s,a) (r + \gamma \max_{a'} Q_*(s',a'))$$

Bellman equation



Policy improvement theorem

$$\pi_0 \rightarrow \pi_1 \rightarrow \pi_2 \rightarrow \dots \rightarrow \pi^*$$

$$\pi' \succ \pi : \forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

Thm. Even if π and π' begin, also

$$\forall s \quad Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s),$$

$$\forall s \quad V_{\pi'}(s) \geq V_{\pi}(s)$$

Proof: $V_{\pi}(s) \leq Q_{\pi}(s, \pi'(s)) = E_{\pi'} [r_{t+1} + \gamma V_{\pi}(s_{t+1}) | s_t = s] \leq$

$$\leq E_{\pi'} [r_{t+1} + \gamma Q_{\pi}(s_{t+1}, \pi'(s_{t+1})) | s_t = s] =$$

$$= E_{\pi'} [r_{t+1} + \gamma r_{t+2} + \gamma^2 V_{\pi}(s_{t+2}) | s_t = s] \leq$$

$$\leq E_{\pi'} [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots | s_t = s] = V_{\pi^*}(s)$$

Policy iteration

$$\pi_{k+1}(s) = \operatorname{argmax}_a Q_{\pi_k}(s, a)$$

